

sPPM and UMT2K on BGL

Gyan Bhanot
IBM Research

LLNL Applications

SPPM – 3d compressible hydrodynamics
Problem size per node fixed, Grid Comm.
Tests weak scaling

UMT2K - unstructured mesh radiation transport
Problem size fixed, domain decomposition,
Node communicates to many neighbors.
Tests strong scaling

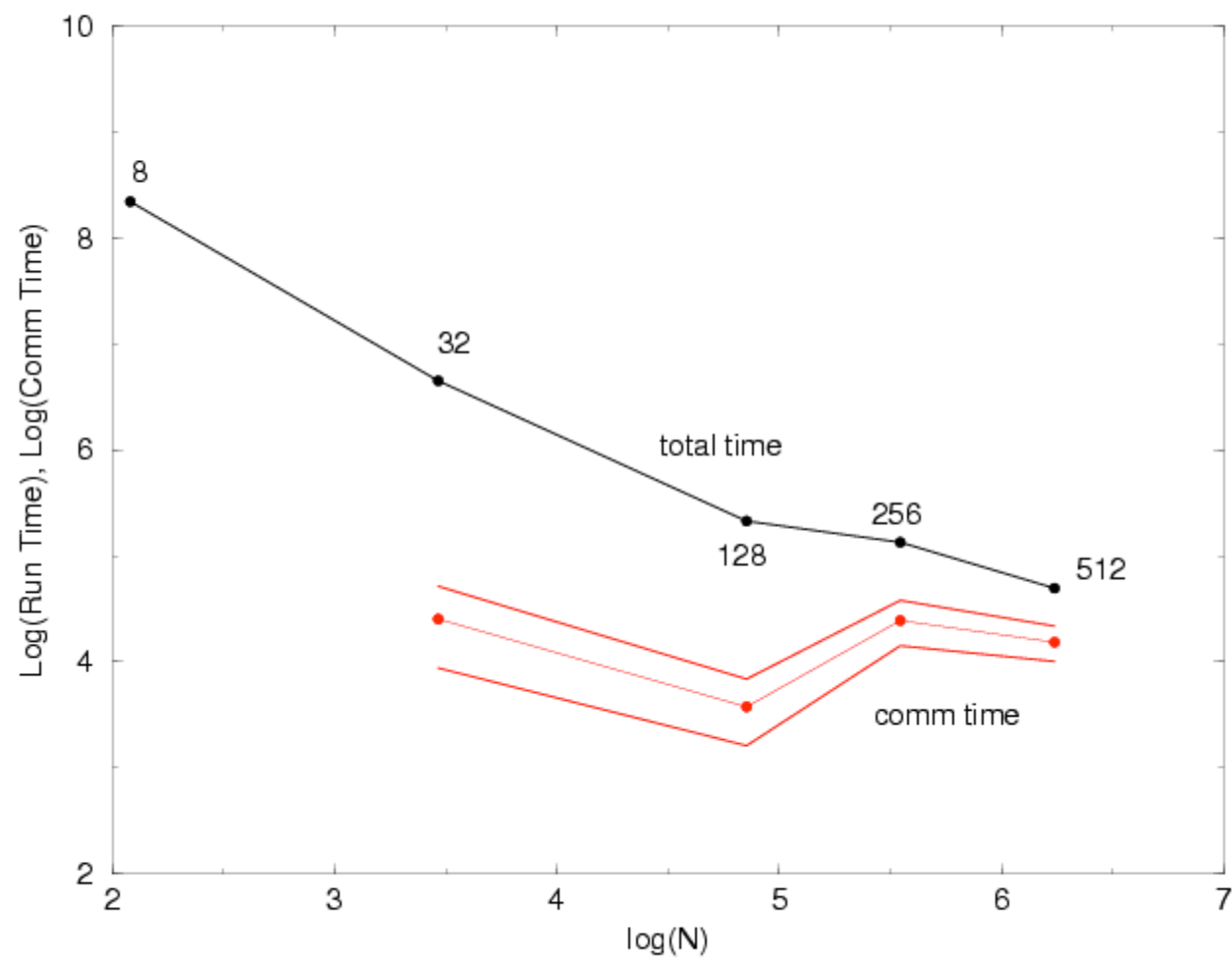
Ran Optimized routines using MPI and blrts_xlc, blrts_xlf

Tested on 8, 32, 128, 256, 512 way

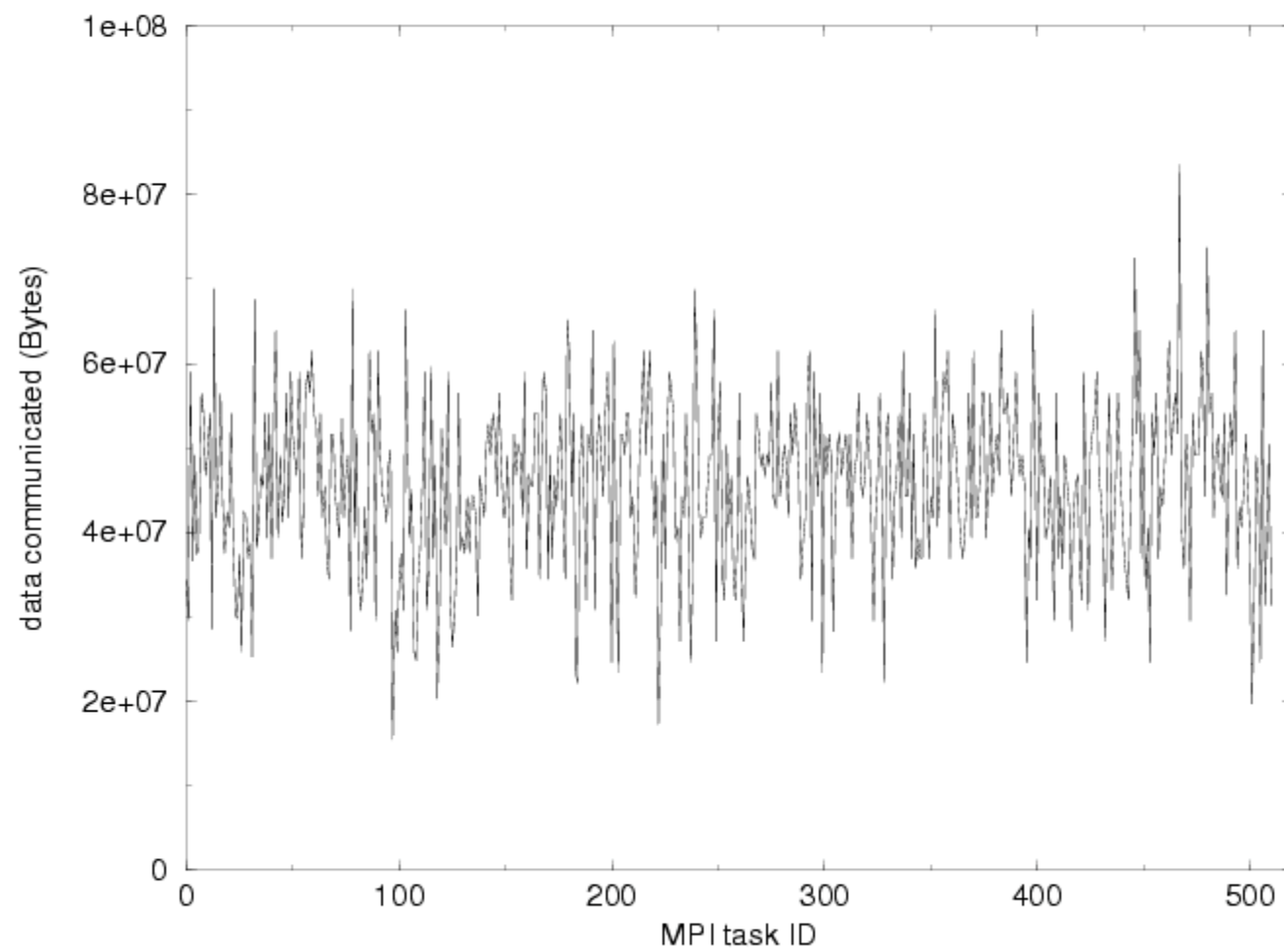
UMT2K Timings (10-12-03)

Nodes	Elapsed (seconds)	Commun (seconds)	Compute (seconds)
512	109	66 +/- 11	43 +/- 11
256	168	80 +/- 17	88 +/- 17
128	206	36 +/- 11	170 +/- 11
32	772	82 +/- 30	691 +/- 30
8	4200		

UMT2K Timings on BGL (Oct 12, 2003)



UMT2K communication on 512 nodes



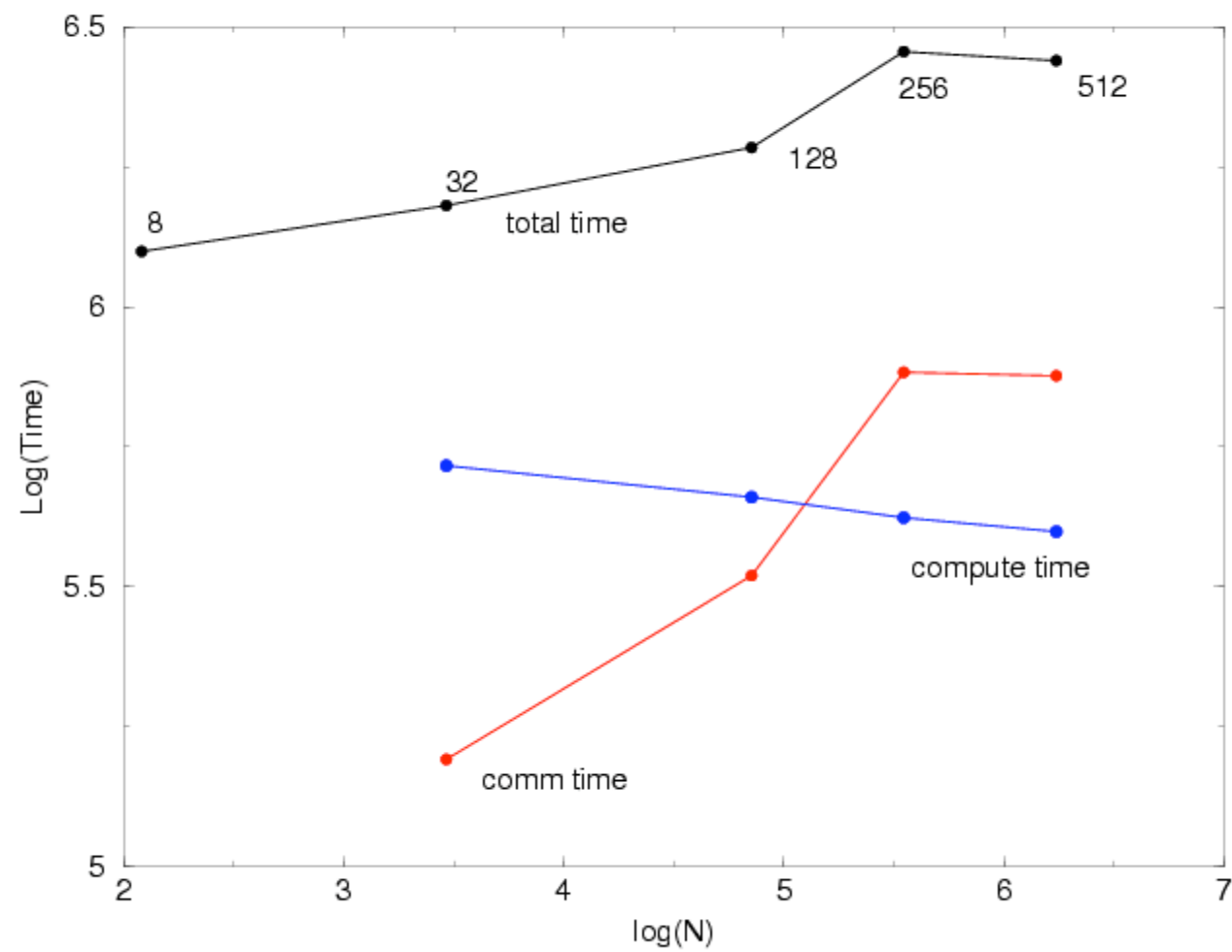
Whats Going on in UMT2K?

- For > 512 nodes test problem is too small. There is not enough work to do per node.
- Problem made worse by load Imbalance. Nodes communicate different amounts of data to different numbers of nodes
- May improve using improved MPI task remapping techniques.

sPPM timings (10-12-03)

Nodes	Elapsed (seconds)	Commun. (seconds)	Compute (seconds)
512	626	357 +/- 15	269 +/- 15
256	636	359 +/- 16	277 +/- 16
128	537	250 +/- 17	287 +/- 17
32	484	180 +/- 11	304 +/- 11
8	445		

SPPM Timings (40 iters) on Oct 12, 2003



What is Going On in sPPM?

- Currently, system uses only single CPU for compute and communication
- MPI_ISEND and MPI_IRECV just return
- All communication is done at MPI_WAIT
- sPPM is written to do computation after MPI_IRECV, MPI_ISEND are posted to overlap computation and communication
- Results in lots of time at MPI_WAIT

SPPM timing on 512 nodes (Oct 12, 2003)

